# What Do Hebbian Learners Learn?

## Reduction Axioms for Iterated Hebbian Learning

**Caleb Schultz Kisby**,

with Saúl Blanco, Larry Moss
Indiana University

**AAAI 2024**
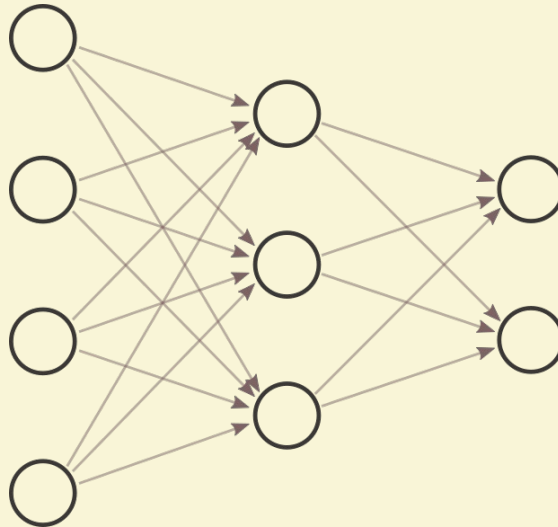February 22, 2024

# Foundations for Neuro-Symbolic AI

From van Harmelen (2022):

> "What are the possible interactions between knowledge and learning? Can reasoning be used as a symbolic prior for learning . . . Can symbolic constraints be enforced on data-driven systems to make them safer? Or less biased? Or can, vice versa, learning be used to yield symbolic knowledge? And if so, how to manage the inherent uncertainty that comes with such learned knowledge . . ."

> ". . . **neuro-symbolic systems currently lack a theory that even begins to ask these questions, let alone answer them.**"

van Harmelen, F. "Preface: The 3rd AI Wave Is Coming, and It Needs a Theory". In: Neuro-Symbolic Artificial Intelligence. Ed. by P. Hitzler and M. Sarker. IOS Press BV, 2022.
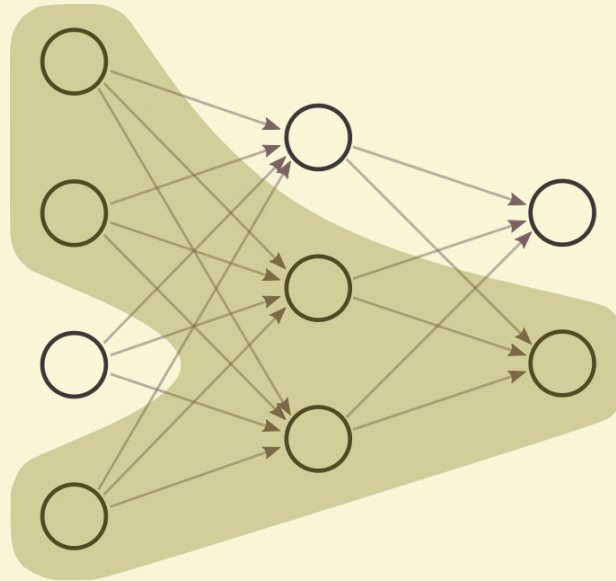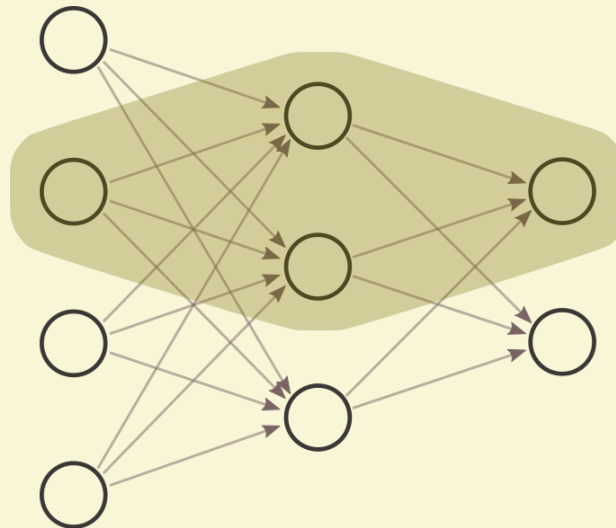
# Neural Network Semantics
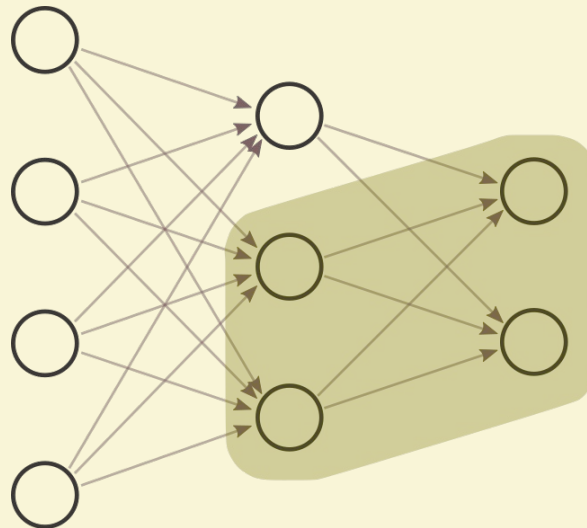
- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net's states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes! Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics
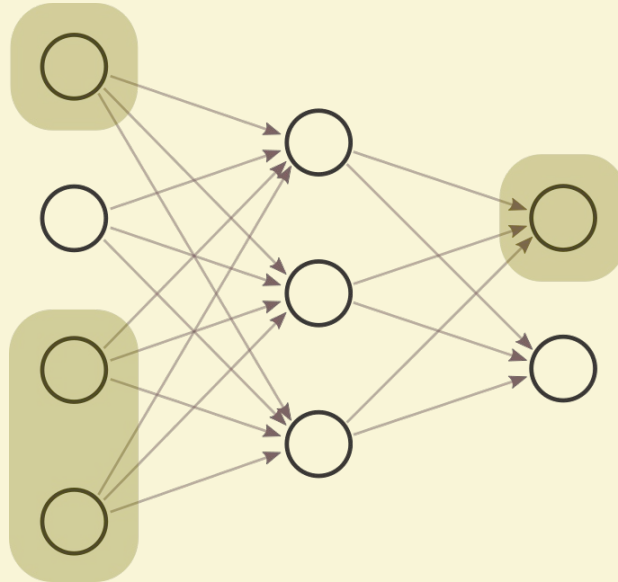
- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net' s states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes!  Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics
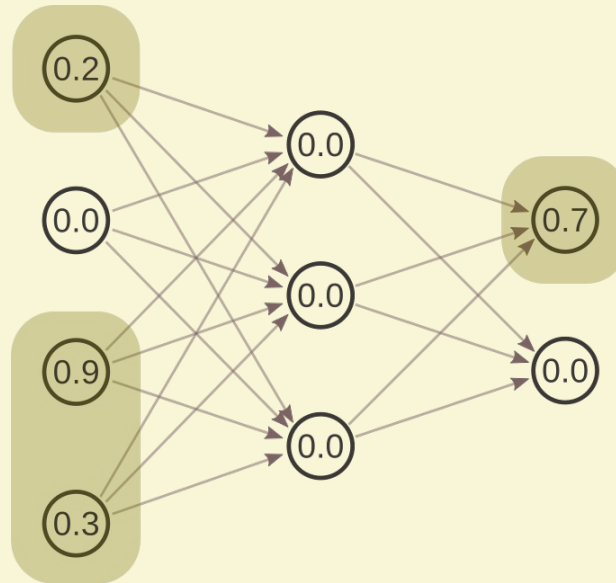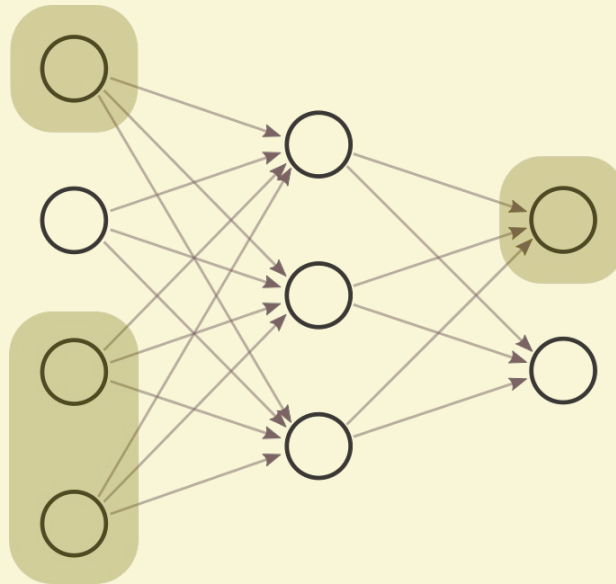
- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net' s states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes!  Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics
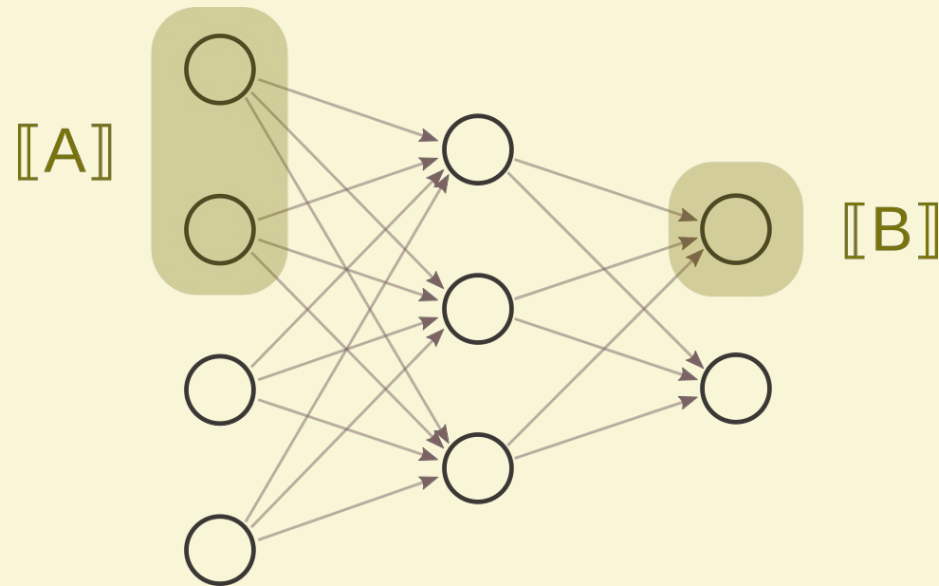
- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net's states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes! Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.
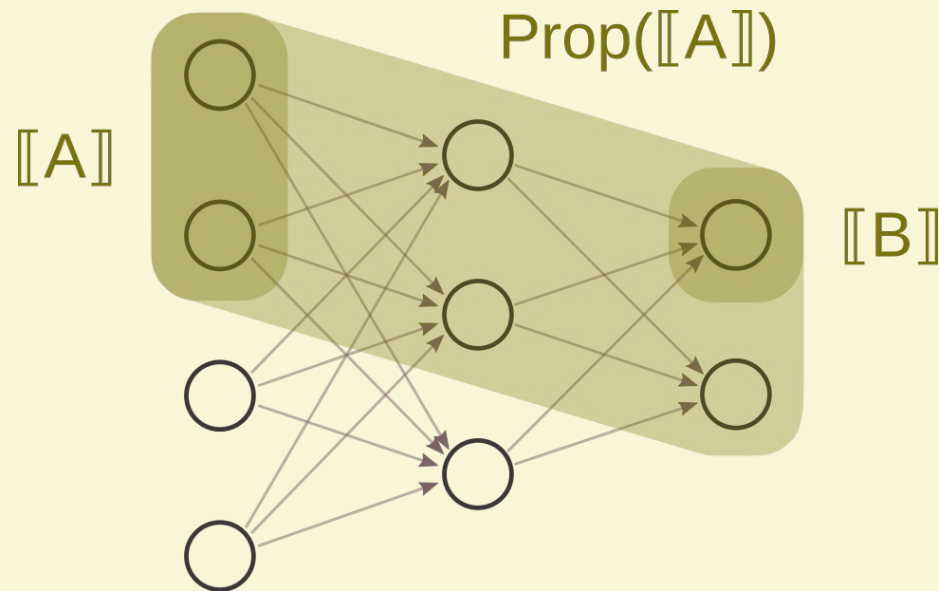
Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics

- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net' s states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes!  Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics

- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net' s states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes! Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics

- **We assume:** The network is weighted, feed-forward, fully-connected, with binary activations. The net's states (activation patterns) are just given by sets of nodes.



- **Key Idea:** Neural networks are not merely black boxes! Instead, think of nets as a kind of (logical) model; The dynamics of its states contain information about its conditional beliefs.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics (Contd.)

- An input state will activate new nodes, which subsequently activate more nodes. The forward propagation **Prop(S)** is the set of all neurons that are eventually activated by S.



The net satisfies A ⇒ B iff Prop(⟦A⟧) ⊇ ⟦B⟧

In other words, the net *classifies A as B*.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Neural Network Semantics (Contd.)

- An input state will activate new nodes, which subsequently activate more nodes. The forward propagation **Prop(S)** is the set of all neurons that are eventually activated by S.



The net satisfies A ⇒ B iff Prop(〚A〛) ⊇ 〚B〛

In other words, the net *classifies A as B*.

Balkenius, C. and Gardenfors, P. Nonmonotonic inferences in neural networks. In KR, 32–39. Morgan Kaufmann, 1991.

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Example: Building a Neural Network

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

⤳⇥

⟦bird⟧

⟦flies⟧

-100

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

⤳

⟦bird⟧

Prop(⟦bird⟧)

⟦flies⟧

-100

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

⇝

⟦bird⟧

⟦flies⟧

-100

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

⤳

⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

# Example: Building a Neural Network



penguin → bird
bird ⇒ flies
¬ (penguin ⇒ flies)

⤳

⟦bird⟧

⟦penguin⟧

-100

⟦flies⟧

Prop(⟦penguin⟧)

# Soundness and Completeness

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Soundness and Completeness

## Soundness

$$\Gamma \vdash A \text{ implies } \Gamma \models A$$

- **Not:** An explanation of a *particular* neural network's behavior

- **But instead:** Sound rules give *high-level* properties for *all* neural networks (of a certain architecture)

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Soundness and Completeness

| Soundness | Completeness |
|---|---|
| $\Gamma \vdash A$ implies $\Gamma \models A$ | $\Gamma \models A$ implies $\Gamma \vdash A$ |

**Soundness**

$\Gamma \vdash A$ implies $\Gamma \models A$

- **Not:** An explanation of a *particular* neural network's behavior

- **But instead:** Sound rules give *high-level* properties for *all* neural networks (of a certain architecture)

**Completeness**

$\Gamma \models A$ implies $\Gamma \vdash A$

- **Equivalently:** Can we build a neural network satisfying the set $\Gamma$ of constraints?

penguin $\rightarrow$ bird
bird $\Rightarrow$ flies
$\neg$ (penguin $\Rightarrow$ flies)

⟦bird⟧

⟦penguin⟧

-100

⟦flies⟧

Prop(⟦penguin⟧)

Leitgeb, H. Neural Network Models of Conditionals. In Introduction to Formal Philosophy, 147–176. Springer, 2018.

# Iterated Hebbian Learning

*Neurons that fire together wire together*



Repeat this update until a fixed point!
i.e. until the weights are "maximally high"

We call the resulting net **Hebb\*(N, ⟦S⟧)**

D. Hebb. The Organization of Behavior. Psychology Press, 1949.

# Iterated Hebbian Learning

*Neurons that fire together wire together*



Repeat this update until a fixed point!
i.e. until the weights are "maximally high"

We call the resulting net **Hebb*(N, ⟦S⟧)**

D. Hebb. The Organization of Behavior. Psychology Press, 1949.

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

⟦puffin⟧

-100

26

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

⟦puffin⟧

pos

Prop(⟦puffin⟧)

27

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

Hebb*(N, ⟦puffin⟧)

⟦puffin⟧

pos

Prop(⟦puffin⟧)

Louis Agassiz Fuertes. Atlantic Puffin (1932). Watercolor and pencil on paper. From *Portraits of New England Birds*. Commonwealth Of Massachusetts, 1932. Edited by Sabrina Schultz Kisby.

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

Hebb*(N, ⟦puffin⟧)

⟦bird⟧

⟦penguin⟧

⟦flies⟧

pos

⟦puffin⟧

pos

Prop(⟦puffin⟧)

29

# Example: Learning Wrecks the Model!



⟦bird⟧

⟦penguin⟧

⟦flies⟧

-100

Prop(⟦penguin⟧)

Hebb*(N, ⟦puffin⟧)

⟦bird⟧

⟦penguin⟧

⟦flies⟧

pos

Prop(⟦penguin⟧)

⟦puffin⟧

pos

Prop(⟦puffin⟧)

# Logic & Formal Semantics

**Syntax.** We consider the language:

$$A, B \in p \mid \neg A \mid A \wedge B \mid \mathbf{K} A \mid \mathbf{T} A$$

We define the duals $\langle \mathbf{K} \rangle, \langle \mathbf{T} \rangle$ as usual. We can express $A \Rightarrow B$ as $\mathbf{T} A \rightarrow B$ ("the typical $A$ is $B$").

**Semantics.** We map each formula to a state:

$$[\![p]\!] = V(p) \quad [\![\neg A]\!] = [\![A]\!]^{\complement} \quad [\![A \wedge B]\!] = [\![A]\!] \cap [\![B]\!]$$

$$[\![\mathbf{K} A]\!] = \{n \mid n \text{ is graph-reachable from } A\}$$
$$[\![\mathbf{T} A]\!] = \text{Prop}([\![A]\!])$$

**Definition.** $N, w \vDash A$ iff $w \in [\![A]\!]$

$$[\![[A]B]\!]_N = [\![B]\!]_{\text{Hebb}^*(N, [\![A]\!])}$$

Can we completely characterize $[A]$'s effect on the net?

# Main Results

**Theorem.**    The following axioms are sound:

$$[\varphi]p \quad\quad\quad\leftrightarrow\quad p \quad\quad\quad \text{for propositions } p$$
$$[\varphi]\neg\psi \quad\quad\leftrightarrow\quad \neg[\varphi]\psi$$
$$[\varphi](\psi \wedge \rho) \quad\leftrightarrow\quad [\varphi]\psi \wedge [\varphi]\rho$$
$$[\varphi]\mathbf{K}\psi \quad\quad\leftrightarrow\quad \mathbf{K}[\varphi]\psi$$
$$[\varphi]\mathbf{T}\psi \quad\quad\leftrightarrow\quad \mathbf{T}([\varphi]\psi \wedge (\mathbf{T}\varphi \vee \mathbf{K}(\mathbf{T}\varphi \vee \mathbf{T}[\varphi]\psi)))$$

**Theorem.**    **Assuming** model building for the base language: For all consistent $\Gamma \subseteq \mathcal{L}$ there is a net $\mathcal{N}$ such that $\mathcal{N} \models \Gamma$.

**Theorem.**    **Assuming** completeness for the base language: $[\varphi]$ is completely axiomatized by the reduction axioms from before.

# Future Work

- **Stabilized Learning Policies.** Hebb* increases the weights until they're maximally large. But stabilized hebbian learning (e.g. Oja's rule) increases weights towards a convergent point.

- **Single-step Update.** We often want guarantees for what the neural network learns *at each step*. This is the heart of AI Alignment.

- **What about Backpropagation?** This is our major long-term goal!

# Check Out Our Poster + Paper!



**Contact:**
Caleb Schultz Kisby
cckisby@iu.edu
https://ais-climber.github.io/